

Supplemental Information

Multivariate analysis of heritable traits

Christoph Lippert^{*,†}, Francesco Paolo Casale^{*}, Barbara Rakitsch, Oliver Stegle^{*,†}

^{*}These authors have contributed equally.

[†]Please address correspondence to lippert@microsoft.com and oliver.stegle@ebi.ac.uk.

Contents

1 LIMIX methods and implementation	1
1.1 Kernel methods for defining composite covariance functions	1
1.1.1 Base covariance models	1
1.1.2 Composite covariance functions	2
1.2 Inference and parameter estimation	3
1.2.1 Maximizing the model likelihood	3
1.2.2 Regularization	3
1.3 Kronecker product identities	4
1.4 Efficient inference for matrix variate mixed models	4
1.4.1 Log likelihood evaluation	5
1.4.2 Estimation of the covariance parameters	6
1.4.3 Estimation of the fixed effects	7
1.4.4 Predictions	9
1.5 Implementation details	9
2 Multivariate linear mixed models for statistical genetics	11
2.1 Introduction	11
2.1.1 Some recently proposed MTMM	11
2.2 Variance Decomposition	12
2.2.1 Single-trait variance decomposition	12
2.2.2 Multiple-trait variance decomposition	12
2.2.3 Phenotype Predictions	13
2.3 Genome-wide Association Studies	13
2.3.1 Univariate GWAS	13
2.3.2 Multivariate GWAS	14
2.3.3 Multi-locus GWAS	15
2.4 PANAMA	15

1 LIMIX methods and implementation

LIMIX is a flexible mixed model framework, allowing to efficiently formulate a variety of different mixed model analyses. Complex covariance functions to define random effects can be combined from simpler building blocks, using addition and multiplication as basic operations (Section 1.1). Inference in these models can be performed by gradient-based optimisation of the marginal likelihood or a regularized version thereof (Section 1.2). Section 1.4 describes how the identical framework can be extended to perform efficient inference in the special case of matrix variate data and Section 1.5 outlines the software implementation and exemplifies the basic usage.

In the most general formulation, we assume that the data \mathbf{y} is multivariate normal distributed with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$

$$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_\theta; \boldsymbol{\Sigma}_\theta), \quad (1.1)$$

where $\boldsymbol{\theta}$ is the set of model parameters. In linear mixed models $\boldsymbol{\mu}_\theta$ typically is a linear function $\boldsymbol{\mu}_\theta = \mathbf{X}\boldsymbol{\beta}$, parametrized by a fixed design matrix \mathbf{X} and fixed-effect parameters $\boldsymbol{\beta}$. The covariance matrix $\boldsymbol{\Sigma}_\theta$ is parametrized by individual variance components or covariance functions, corresponding to random effects in the model. In this instance, the set of model parameters $\boldsymbol{\theta}$ is given by the union of the fixed effects $\boldsymbol{\beta}$ and the variance parameters.

1.1 Kernel methods for defining composite covariance functions

Defining property of a valid covariance matrix is a positive semi-definite symmetric matrix. In LIMIX we utilize a key insight from the kernel machines literature that allows to construct complex covariance matrices by combining a set of basic covariance models. First, we define a set of base covariance functions used in LIMIX (Section 1.1.1). In Section 1.1.2, we then define basic operations to combine two base covariance matrices \mathbf{A} and \mathbf{B} using sum and multiplication operators, which obey positive semi-definiteness of the resulting matrix \mathbf{C} . By iteratively applying these operations, it is possible to build complex covariance matrices, thereby allowing for extensive modeling flexibility.

1.1.1 Base covariance models

LIMIX provides implementations for a range of covariance matrices, where the most important examples are defined here. A more rigorous treatment of covariance functions including the definition of more extensive examples can be found in [1]. Each covariance is parameterised by a set of parameters $\boldsymbol{\theta}$, giving rise to a semi-positive definite covariance matrix $\mathbf{C}(\boldsymbol{\theta})$. For inference, the matrix derivatives of \mathbf{C} w.r.t θ_i are needed, see also Section 1.2.

Fixed covariance matrix

If the covariance \mathbf{A} between the samples is observed, inference requires to merely learn a single amplitude parameter $\sigma_g^2 > 0$:

$$\mathbf{C}(\boldsymbol{\theta}) = \sigma_g^2 \mathbf{A}, \quad \frac{\partial \mathbf{C}}{\partial \sigma_g^2} = \mathbf{A} \quad (1.2)$$

Here, $\boldsymbol{\theta} = [\sigma_g^2]$.

Linear covariance function

The linear covariance function can be derived by marginalizing out the weights in a random effect model $\mathbf{X}\boldsymbol{\beta}$, assuming $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \sigma_g^2 \mathbf{I})$. This results in a covariance of the form $\mathbf{C}(\boldsymbol{\theta}) = \sigma_g^2 \mathbf{X} \mathbf{X}^\top$, where again $\boldsymbol{\theta} = [\sigma_g^2]$. In genetics, this covariance can be used to model polygenic effects if \mathbf{X} corresponds to genome-wide SNP data, e.g. [2, 3].

Free form covariance model

In practice, the covariance between phenotypes is not observed and hence its parameters need to be inferred from data. The most general form is the free form covariance, which is a general $P \times P$ matrix that obeys the semi positive-definiteness constraints. To this end, we parametrize the covariance using cholesky factors, as any semi-positive definite matrix has a valid cholesky decomposition.

$$\mathbf{C}(\boldsymbol{\theta}) = \mathbf{L} \mathbf{L}^\top, \quad \frac{\partial \mathbf{C}}{\partial \mathbf{L}} = \mathbf{L} \partial(\mathbf{L})^\top + \partial(\mathbf{L}) \mathbf{L}^\top \quad (1.3)$$

where \mathbf{L} is a lower triangular matrix of covariance parameters, i.e. $\boldsymbol{\theta} = [L_{00}, L_{10}, L_{11}, \dots, L_{0n}, L_{nn}]$. The gradient with respect to a matrix field L_{ij} can be obtained by applying the chain rule $\frac{\partial \mathbf{C}}{\partial L_{ij}} = \frac{\partial \mathbf{C}}{\partial \mathbf{L}} \frac{\partial \mathbf{L}}{\partial L_{ij}}$.

Low-rank covariance model

We can restrict the phenotypes to lie in a lower-dimensional linear subspace, by using a low-rank parameterization of the covariance matrix

$$\mathbf{C}(\boldsymbol{\theta}) = \mathbf{Z} \mathbf{Z}^\top, \quad \frac{\partial \mathbf{C}}{\partial \mathbf{Z}} = 2\mathbf{Z} \quad (1.4)$$

where \mathbf{Z} is a $P \times K$ matrix, with $K < P$, and the parameters are its entries, i.e. $\boldsymbol{\theta} = [Z_{00}, Z_{01}, \dots, Z_{0k}, \dots, Z_{n0}, \dots, Z_{nk}]$. The gradient with respect to an entry Z_{ij} can then again be obtained by using the chain rule $\frac{\partial \mathbf{C}}{\partial Z_{ij}} = \frac{\partial \mathbf{C}}{\partial \mathbf{Z}} \frac{\partial \mathbf{Z}}{\partial Z_{ij}}$.

Diagonal covariance matrix

The diagonal covariance matrix is often used to model the noise between the phenotypes. By restricting the matrix to be diagonal, we implicitly assume that the noise is independently, but not identically distributed

$$\mathbf{C}(\boldsymbol{\theta}) = \text{diag}(\mathbf{d}), \quad \frac{\partial \mathbf{C}}{\partial \mathbf{d}} = \mathbf{I}, \quad (1.5)$$

where $\boldsymbol{\theta} = [d_1, \dots, d_n]$.

1.1.2 Composite covariance functions

First, the sum of two covariance matrices is a covariance matrix

$$\mathbf{A} + \mathbf{B} = \mathbf{C}. \quad (1.6)$$

Additive combinations of covariance matrices allow for variance decomposition and can be interpreted as multiple independent random effects. Second, the point-wise product of two covariance matrices is a valid covariance matrix,

$$\mathbf{A} \odot \mathbf{B} = \mathbf{C} \quad (1.7)$$

and the Kronecker product of two covariances matrices is as well (see also Section 1.4)

$$\mathbf{A} \otimes \mathbf{B} = \mathbf{C}. \quad (1.8)$$

The gradient of the composite covariance functions can be obtained by applying iteratively the sum and product rule

$$\frac{\partial}{\partial \boldsymbol{\theta}} (\mathbf{A} + \mathbf{B}) = \frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{A} + \frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{B} \quad (1.9)$$

$$\frac{\partial}{\partial \boldsymbol{\theta}} (\mathbf{A} \odot \mathbf{B}) = \frac{\partial}{\partial \boldsymbol{\theta}} (\mathbf{A}) \odot \mathbf{B} + \mathbf{A} \odot \frac{\partial}{\partial \boldsymbol{\theta}} (\mathbf{B}) \quad (1.10)$$

$$\frac{\partial}{\partial \boldsymbol{\theta}} (\mathbf{A} \otimes \mathbf{B}) = \frac{\partial}{\partial \boldsymbol{\theta}} (\mathbf{A}) \otimes \mathbf{B} + \mathbf{A} \otimes \frac{\partial}{\partial \boldsymbol{\theta}} (\mathbf{B}). \quad (1.11)$$

1.2 Inference and parameter estimation

1.2.1 Maximizing the model likelihood

We perform parameter estimation by maximizing the log likelihood. The gradient of the log likelihood is optimized using LBFGS [4].

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\theta}_i} \log \mathcal{N}(\mathbf{y} | \boldsymbol{\mu}_{\boldsymbol{\theta}}; \boldsymbol{\Sigma}_{\boldsymbol{\theta}}) = & -\frac{1}{2} \text{tr} \left(\boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \frac{\partial \boldsymbol{\Sigma}_{\boldsymbol{\theta}}}{\partial \boldsymbol{\theta}_i} \right) + \frac{1}{2} (\mathbf{y} - \boldsymbol{\mu}_{\boldsymbol{\theta}})^{\top} \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \frac{\partial \boldsymbol{\Sigma}_{\boldsymbol{\theta}}}{\partial \boldsymbol{\theta}_i} \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} (\mathbf{y} - \boldsymbol{\mu}_{\boldsymbol{\theta}}) \\ & - (\mathbf{y} - \boldsymbol{\mu}_{\boldsymbol{\theta}})^{\top} \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \frac{\partial (\mathbf{y} - \boldsymbol{\mu}_{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}_i} \end{aligned} \quad (1.12)$$

Repeat optimizations from independent starting points and informative initializations are employed to improve convergence of the optimization procedure.

1.2.2 Regularization

Estimating the covariance matrix between phenotypes is hindered by the large parameter space: for free-form covariance matrices, the number of parameters is growing quadratically with the number of phenotypes, while the number of new data points is only growing linearly. If not accounted for, this will lead to overfitting, i.e. the model does not only explain the signal of the data but also the noise, which in turn leads to a bad generalization behavior and a loss of interpretability. A natural way to prevent this is to add a regularization term over the covariance matrix

$$\min_{\boldsymbol{\theta}} \underbrace{-\log \mathcal{N}(\mathbf{y} | \boldsymbol{\mu}_{\boldsymbol{\theta}}; \boldsymbol{\Sigma}_{\boldsymbol{\theta}})}_{\text{data fit}} + \underbrace{\text{Reg}(\boldsymbol{\Sigma}_{\boldsymbol{\theta}})}_{\text{regularizer}}. \quad (1.13)$$

While the first term maximizes the fit between the data and the model, the second term acts as a penalizer that prevents the model from becoming too complex. LIMIX makes available a sum-of-squares penalty on the off-diagonal entries of the covariance matrix $\text{Reg}(\boldsymbol{\Sigma}_{\boldsymbol{\theta}}) = \lambda \sum_{i \neq j} (\boldsymbol{\Sigma}_{\boldsymbol{\theta}})_{ij}^2$, or on its inverse $\text{Reg}(\boldsymbol{\Sigma}_{\boldsymbol{\theta}}) = \lambda \sum_{i \neq j} (\boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1})_{ij}^2$, where λ is the trade-off parameter between fitting the data and regularizing. In multi-trait models, in particular for larger numbers of traits (eQTL analysis in yeast) we considered the latter. Both this penalizations do not penalize the diagonal elements, to retain an unbiased estimate of the marginal heritabilities. From the Bayesian perspective, the regularization term is equivalent to an isotropic Gaussian prior, with zero mean.

Out of sample prediction In practice, it is often hard to find the right trade-off between model fit and regularization. We used out of sample prediction accuracy to find the appropriate λ . Predictions for new data points can be carried out by conditioning on the observed data

$$\begin{aligned} p(\mathbf{y}^* | \mathbf{y}, \boldsymbol{\theta}) = & \mathcal{N}(\mathbf{y}^* | \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{X}^*) + \boldsymbol{\Sigma}_{\boldsymbol{\theta}}(\mathbf{X}^*, \mathbf{X}) \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} (\mathbf{y} - \boldsymbol{\mu}_{\boldsymbol{\theta}}); \\ & \boldsymbol{\Sigma}_{\boldsymbol{\theta}}(\mathbf{X}^*, \mathbf{X}^*) - \boldsymbol{\Sigma}_{\boldsymbol{\theta}}(\mathbf{X}^*, \mathbf{X}) \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}}(\mathbf{X}, \mathbf{X}^*)), \end{aligned} \quad (1.14)$$

where $\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{X}^*)$ is the mean function on the test inputs, $\boldsymbol{\Sigma}_{\boldsymbol{\theta}}(\mathbf{X}, \mathbf{X})^*$ depicts the covariance between the training and the test samples, and $\boldsymbol{\Sigma}_{\boldsymbol{\theta}}(\mathbf{X}^*, \mathbf{X}^*)$ between the test samples.

1.3 Kronecker product identities

Let \mathbf{A} be a $N \times M$ matrix, and \mathbf{B} be a $P \times Q$ matrix. The Kronecker product $\mathbf{A} \otimes \mathbf{B}$ is a $NP \times MQ$ matrix and defined as follows

$$\mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} A_{11}\mathbf{B} & \dots & A_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ A_{m1}\mathbf{B} & \dots & A_{mn}\mathbf{B} \end{pmatrix} \quad (1.15)$$

Due to its block structure, the Kronecker product has a number of nice properties that speed up inference[5, 6]. In particular, the dot product between two Kronecker products is a Kronecker product again

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AC} \otimes \mathbf{BD}, \quad (1.16)$$

where $\mathbf{C} \in \mathbb{R}^{M \times R}$ and $\mathbf{D} \in \mathbb{R}^{Q \times S}$. It can be evaluated in $O(NMR + PQS)$ time, while the naive runtime is $O(NPMQRS)$.

Let vec be an operation that concatenates the columns of a $N \times M$ matrix into a vector of length $N \cdot M$. The product of Kronecker product and a vectorized matrix can be computed efficiently

$$(\mathbf{A} \otimes \mathbf{B}) \text{vec}(\mathbf{Y}) = \text{vec}(\mathbf{BYA}^\top) \quad (1.17)$$

bringing the runtime down to $O(\min(PQM + PMN, QMN + PQN))$ from $O(NPMQ)$.

Let $\mathbf{U}_A \mathbf{S}_A \mathbf{U}_A^\top$ be the eigenvalue decomposition of $\mathbf{A} \in \mathbb{R}^{N \times N}$ and $\mathbf{U}_B \mathbf{S}_B \mathbf{U}_B^\top$ the eigenvalue decomposition of $\mathbf{B} \in \mathbb{R}^{P \times P}$. The eigenvalue decomposition of the Kronecker product $\mathbf{A} \otimes \mathbf{B}$ can be synthesized by the eigenvalue decompositions of its components

$$\mathbf{A} \otimes \mathbf{B} = (\mathbf{U}_A \otimes \mathbf{U}_B)(\mathbf{S}_A \otimes \mathbf{S}_B)(\mathbf{U}_A^\top \otimes \mathbf{U}_B^\top) \quad (1.18)$$

leading to a runtime reduction from $O(N^3P^3)$ to $O(N^3 + P^3)$.

1.4 Efficient inference for matrix variate mixed models

In Section 1.1, we gave a general overview over existing covariance functions. We now consider the special case in which the covariance matrix can be written as a sum of two Kronecker products and the fixed design matrix also exhibits Kronecker structure

$$\mathcal{N} \left(\text{vec}(\mathbf{Y}) \mid \sum_{j=1}^J (\mathbf{A}_j \otimes \mathbf{X}_j) \text{vec}(\mathbf{B}_j); \underbrace{\mathbf{C}(\boldsymbol{\theta}_C) \otimes \mathbf{R}(\boldsymbol{\theta}_R)}_{\text{signal}} + \underbrace{\boldsymbol{\Sigma}(\boldsymbol{\theta}_\Sigma) \otimes \mathbf{I}}_{\text{noise}} \right), \quad (1.19)$$

Here, \mathbf{Y} denotes the $N \times P$ data matrix with N samples and P phenotypes. For the fixed effect j , the matrix $\mathbf{B}_j \in \mathbb{R}^{M_j \times D_j}$ denotes the fixed-effects, $\mathbf{A}_j \in \mathbb{R}^{P \times M_j}$ is the design matrix of the column effects and $\mathbf{X}_j \in \mathbb{R}^{N \times D_j}$ of the row effects. We define $\mathbf{R}(\boldsymbol{\theta}_R)$ as the signal row covariance of the data matrix and $\mathbf{C}(\boldsymbol{\theta}_C)$ as the signal column covariance matrix. Similarly, the noise column covariance matrix is given by $\boldsymbol{\Sigma}(\boldsymbol{\theta}_\Sigma)$. For the sake of clarity, we further assume that the noise row covariance matrix is the identity \mathbf{I} , the extension to an arbitrary noise row covariance matrix is straightforward. For notational convenience, we will also drop the dependence of the covariance matrices on additional hyperparameters from now on.

The merits of this class of models lie in its tractability: as we show in the following, efficient inference can be performed in $O(N^3 + P^3)$ time and in $O(N^2 + P^2)$ space, whereas an arbitrary covariance matrix of size NP has runtime $O(N^3P^3)$ and a memory requirement of $O(N^2P^2)$.

1.4.1 Log likelihood evaluation

The log likelihood of the data is given by:

$$\log \mathcal{L} = -\frac{NP}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{C} \otimes \mathbf{R} + \boldsymbol{\Sigma} \otimes \mathbf{I}| - \frac{1}{2} \text{vec}(\mathbf{Y}_r)^\top (\mathbf{C} \otimes \mathbf{R} + \boldsymbol{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}_r), \quad (1.20)$$

where \mathbf{Y}_r is the residual phenotype, after the fixed effects have been subtracted from the data:

$$\begin{aligned} \text{vec}(\mathbf{Y}_r) &= \text{vec}(\mathbf{Y}) - \sum_{j=1}^J (\mathbf{A}_j \otimes \mathbf{X}_j) \text{vec}(\mathbf{B}_j) \\ &= \text{vec} \left(\mathbf{Y} - \sum_{j=1}^J \mathbf{X}_j \mathbf{B}_j \mathbf{A}_j^\top \right) \end{aligned} \quad (1.21)$$

Let $\mathbf{U}_\Sigma \mathbf{S}_\Sigma \mathbf{U}_\Sigma^\top$ be the eigenvalue decomposition of $\boldsymbol{\Sigma}$. As shown in [7, 3], we can whiten the covariance matrix in a two-step procedure. First, we whiten the noise by factoring it out. Second, we exploit the fact that the eigenvalue decomposition of a Kronecker product is compatible with a constant diagonal matrix to whiten the remaining covariance matrix:

$$\begin{aligned} \mathbf{K} &= \mathbf{C} \otimes \mathbf{R} + \boldsymbol{\Sigma} \otimes \mathbf{I} \\ &\stackrel{(1.18)}{=} \mathbf{C} \otimes \mathbf{R} + \mathbf{U}_\Sigma \mathbf{S}_\Sigma \mathbf{U}_\Sigma^\top \otimes \mathbf{I} \\ &\stackrel{(1.16)}{=} \left(\mathbf{U}_\Sigma \mathbf{S}_\Sigma^{\frac{1}{2}} \otimes \mathbf{I} \right) \left(\underbrace{\mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_\Sigma^\top \mathbf{C} \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \otimes \mathbf{R} + \mathbf{I} \otimes \mathbf{I}}_{\tilde{\mathbf{C}}} \right) \left(\mathbf{S}_\Sigma^{\frac{1}{2}} \mathbf{U}_\Sigma^\top \otimes \mathbf{I} \right) \\ &\stackrel{(1.18)}{=} \left(\mathbf{U}_\Sigma \mathbf{S}_\Sigma^{\frac{1}{2}} \otimes \mathbf{I} \right) \left(\mathbf{U}_{\tilde{\mathbf{C}}} \mathbf{S}_{\tilde{\mathbf{C}}} \mathbf{U}_{\tilde{\mathbf{C}}}^\top \otimes \mathbf{U}_R \mathbf{S}_R \mathbf{U}_R^\top + \mathbf{I} \otimes \mathbf{I} \right) \left(\mathbf{S}_\Sigma^{\frac{1}{2}} \mathbf{U}_\Sigma^\top \otimes \mathbf{I} \right) \\ &\stackrel{(1.16)}{=} \left(\mathbf{U}_\Sigma \mathbf{S}_\Sigma^{\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \otimes \mathbf{U}_R \right) \left(\mathbf{S}_{\tilde{\mathbf{C}}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I} \right) \left(\mathbf{U}_{\tilde{\mathbf{C}}}^\top \mathbf{S}_\Sigma^{\frac{1}{2}} \mathbf{U}_\Sigma^\top \otimes \mathbf{U}_R^\top \right), \end{aligned} \quad (1.22)$$

where $\mathbf{U}_{\tilde{\mathbf{C}}} \mathbf{S}_{\tilde{\mathbf{C}}} \mathbf{U}_{\tilde{\mathbf{C}}}^\top$ is the eigenvalue decomposition of $\tilde{\mathbf{C}}$. The log determinant can then be written as

$$\log |\mathbf{C} \otimes \mathbf{R} + \boldsymbol{\Sigma} \otimes \mathbf{I}| \stackrel{|AB|=|A|\cdot|B|}{=} \log |\mathbf{S}_\Sigma \otimes \mathbf{I}| + \log |\mathbf{S}_{\tilde{\mathbf{C}}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I}| \quad (1.23)$$

$$\stackrel{|A \otimes B|=|A|^P \cdot |B|^N}{=} N \sum_{p=1}^P \log \mathbf{S}_\Sigma[p, p] + \sum_{p=1}^P \sum_{n=1}^N \log (\mathbf{S}_{\tilde{\mathbf{C}}}[p, p] \mathbf{S}_R[n, n] + 1) \quad (1.24)$$

We can evaluate the squared form efficiently as follows

$$\begin{aligned} &\text{vec}(\mathbf{Y}_r)^\top (\mathbf{C} \otimes \mathbf{R} + \boldsymbol{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}_r) \\ &\stackrel{(1.22)}{=} \text{vec}(\mathbf{Y}_r)^\top \left[\left(\mathbf{U}_\Sigma \mathbf{S}_\Sigma^{\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \otimes \mathbf{U}_R \right) \left(\mathbf{S}_{\tilde{\mathbf{C}}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I} \right) \left(\mathbf{U}_{\tilde{\mathbf{C}}}^\top \mathbf{S}_\Sigma^{\frac{1}{2}} \mathbf{U}_\Sigma^\top \otimes \mathbf{U}_R^\top \right) \right]^{-1} \text{vec}(\mathbf{Y}_r) \\ &= \left[\left(\mathbf{U}_\Sigma \mathbf{S}_\Sigma^{\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \otimes \mathbf{U}_R \right)^{-1} \text{vec}(\mathbf{Y}_r) \right]^\top \left(\mathbf{S}_{\tilde{\mathbf{C}}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I} \right)^{-1} \left[\left(\mathbf{U}_\Sigma \mathbf{S}_\Sigma^{\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \otimes \mathbf{U}_R \right)^{-1} \text{vec}(\mathbf{Y}_r) \right] \\ &= \left[\left(\mathbf{U}_{\tilde{\mathbf{C}}}^\top \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_\Sigma^\top \otimes \mathbf{U}_R^\top \right) \text{vec}(\mathbf{Y}_r) \right]^\top \left(\mathbf{S}_{\tilde{\mathbf{C}}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I} \right)^{-1} \left[\left(\mathbf{U}_{\tilde{\mathbf{C}}}^\top \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_\Sigma^\top \otimes \mathbf{U}_R^\top \right) \text{vec}(\mathbf{Y}_r) \right] \\ &\stackrel{(1.17)}{=} \text{vec} \left[\mathbf{U}_R^\top \mathbf{Y}_r \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \right]^\top \left(\mathbf{S}_{\tilde{\mathbf{C}}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I} \right)^{-1} \text{vec} \left[\mathbf{U}_R^\top \mathbf{Y}_r \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \right] \\ &= \text{vec} \left[\mathbf{U}_R^\top \mathbf{Y}_r \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \right]^\top \text{vec} \left[\mathbf{D} \odot \mathbf{U}_R^\top \mathbf{Y}_r \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \right], \end{aligned} \quad (1.25)$$

where \mathbf{D} is a $N \times P$ matrix defined by the entries

$$\mathbf{D}[n, p] = \frac{1}{\mathbf{S}_{\tilde{\mathbf{C}}}[p, p] \otimes \mathbf{S}_R[n, n] + 1}. \quad (1.26)$$

The inverse of the covariance term times a vector is therewith:

$$(\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}_r) = \left(\mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R \right) \text{vec} \left[\underbrace{\mathbf{D} \odot \mathbf{U}_R^T \mathbf{Y}_r \mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}}}_{\tilde{\mathbf{Y}}_r} \right] \quad (1.27)$$

1.4.2 Estimation of the covariance parameters

We want to evaluate the gradient with respect to a particular covariance parameter $\theta \in \{\boldsymbol{\theta}_C, \boldsymbol{\theta}_R, \boldsymbol{\theta}_{\Sigma}\}$. The derivative consists of the derivative of the determinant term and the derivative of the squared form. Each of these is given separately for the row covariance parameter $\theta_r \in \boldsymbol{\theta}_R$. The gradients for the column covariance parameters can be derived analogously and are omitted for brevity.

We start with the gradient of the log determinant:

$$\begin{aligned} & \frac{\partial}{\partial \theta_r} \log |\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I}| \\ &= \text{tr} \left((\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \left(\mathbf{C} \otimes \frac{\partial}{\partial \theta_r} \mathbf{R} \right) \right) \\ &\stackrel{(1.22)}{=} \text{tr} \left(\left[\left(\mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R \right) (\mathbf{S}_{\tilde{C}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I}) \left(\mathbf{U}_{\tilde{C}}^T \mathbf{S}_{\Sigma}^{\frac{1}{2}} \mathbf{U}_{\Sigma}^T \otimes \mathbf{U}_R^T \right) \right]^{-1} \left(\mathbf{C} \otimes \frac{\partial}{\partial \theta_r} \mathbf{R} \right) \right) \\ &\stackrel{(AB)^{-1} = B^{-1}A^{-1}}{=} \text{tr} \left(\left(\mathbf{U}_{\tilde{C}}^T \mathbf{S}_{\Sigma}^{\frac{1}{2}} \mathbf{U}_{\Sigma}^T \otimes \mathbf{U}_R^T \right)^{-1} (\mathbf{S}_{\tilde{C}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I})^{-1} \left(\mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R \right)^{-1} \left(\mathbf{C} \otimes \frac{\partial}{\partial \theta_r} \mathbf{R} \right) \right) \\ &\stackrel{\text{tr}(AB) = \text{tr}(BA)}{=} \text{tr} \left((\mathbf{S}_{\tilde{C}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I})^{-1} \left(\mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R \right)^{-1} \left(\mathbf{C} \otimes \frac{\partial}{\partial \theta_r} \mathbf{R} \right) \left(\mathbf{U}_{\tilde{C}}^T \mathbf{S}_{\Sigma}^{\frac{1}{2}} \mathbf{U}_{\Sigma}^T \otimes \mathbf{U}_R^T \right)^{-1} \right) \\ &\stackrel{(1.16)}{=} \text{tr} \left((\mathbf{S}_{\tilde{C}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I}) \left(\mathbf{U}_{\tilde{C}}^T \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\Sigma}^T \mathbf{C} \mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R^T \left(\frac{\partial}{\partial \theta_r} \mathbf{R} \right) \mathbf{U}_R \right) \right) \\ &\stackrel{\text{def}(\tilde{\mathbf{C}})}{=} \text{tr} \left((\mathbf{S}_{\tilde{C}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I}) \left(\mathbf{U}_{\tilde{C}}^T \tilde{\mathbf{C}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R^T \left(\frac{\partial}{\partial \theta_r} \mathbf{R} \right) \mathbf{U}_R \right) \right) \\ &\stackrel{\text{tr}(AB) = \text{vec}(A)^T \text{vec}(B)}{=} \text{diag} (\mathbf{S}_{\tilde{C}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I})^T \left[\text{diag} (\mathbf{S}_{\tilde{C}}) \otimes \text{diag} \left(\mathbf{U}_R^T \left(\frac{\partial}{\partial \theta_r} \mathbf{R} \right) \mathbf{U}_R \right) \right] \end{aligned} \quad (1.28)$$

The gradient of the squared form is given by:

$$\begin{aligned} & \frac{\partial}{\partial \theta_r} \text{vec}(\mathbf{Y}_r)^T (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}_r) \\ &= -\text{vec}(\mathbf{Y}_r)^T (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \left(\mathbf{C} \otimes \frac{\partial}{\partial \theta_r} \mathbf{R} \right) (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}_r) \\ &\stackrel{(1.27)}{=} -\text{vec}(\tilde{\mathbf{Y}}_r)^T \left(\mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R \right)^{-1} \left(\mathbf{C} \otimes \frac{\partial}{\partial \theta_r} \mathbf{R} \right) \left(\mathbf{U}_{\tilde{C}}^T \mathbf{S}_{\Sigma}^{\frac{1}{2}} \mathbf{U}_{\Sigma}^T \otimes \mathbf{U}_R^T \right)^{-1} \text{vec}(\tilde{\mathbf{Y}}_r) \\ &\stackrel{(1.16)}{=} -\text{vec}(\tilde{\mathbf{Y}}_r)^T \left[\mathbf{U}_{\tilde{C}}^T \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\Sigma}^T \mathbf{C} \mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R^T \left(\frac{\partial}{\partial \theta_r} \mathbf{R} \right) \mathbf{U}_R \right] \text{vec}(\tilde{\mathbf{Y}}_r) \\ &\stackrel{\text{def}(\tilde{\mathbf{C}})}{=} -\text{vec}(\tilde{\mathbf{Y}}_r)^T \left[\mathbf{S}_{\tilde{C}} \otimes \mathbf{U}_R^T \left(\frac{\partial}{\partial \theta_r} \mathbf{R} \right) \mathbf{U}_R \right] \text{vec}(\tilde{\mathbf{Y}}_r) \\ &\stackrel{(1.17)}{=} -\text{vec}(\tilde{\mathbf{Y}}_r)^T \text{vec} \left[\mathbf{U}_R^T \left(\frac{\partial}{\partial \theta_r} \mathbf{R} \right) \mathbf{U}_R \tilde{\mathbf{Y}}_r \mathbf{S}_{\tilde{C}} \right] \end{aligned} \quad (1.29)$$

1.4.3 Estimation of the fixed effects

Gradient Evaluation We first evaluate the gradient with respect to a single fixed effect

$$\begin{aligned}
& \frac{\partial}{\partial \mathbf{B}_k}_{a,b} \left(-\frac{1}{2} \text{vec}(\mathbf{Y}_r)^\top (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}_r) \right) \\
&= -\text{vec}(\mathbf{Y}_r)^\top (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \text{vec} \left(\frac{\partial \mathbf{Y}_r}{\partial \mathbf{B}_k}_{a,b} \right) \\
&\stackrel{(1.27)}{=} -\text{vec} \left(\tilde{\mathbf{Y}}_r \right)^\top \left(\mathbf{U}_{\tilde{\mathbf{C}}}^\top \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\Sigma}^\top \otimes \mathbf{U}_{\mathbf{R}}^\top \right) \text{vec} \left(\frac{\partial \mathbf{Y}_r}{\partial \mathbf{B}_k}_{a,b} \right) \\
&\stackrel{(1.17)}{=} -\text{vec} \left(\tilde{\mathbf{Y}}_r \right) \text{vec} \left(\mathbf{U}_{\mathbf{R}}^\top \left(\frac{\partial \mathbf{Y}_r}{\partial \mathbf{B}_k}_{a,b} \right) \mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \right) \\
&= -\text{vec} \left(\tilde{\mathbf{Y}}_r \right)^\top \text{vec} \left(\mathbf{U}_{\mathbf{R}}^\top \left([\mathbf{X}_k]_{:,a} [\mathbf{A}_k]_{:,b}^\top \right) \mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \right) \\
&\stackrel{\text{tr}(AB) = \text{vec}(A)^\top \text{vec}(B)}{=} -\text{tr} \left(\tilde{\mathbf{Y}}_r^\top \mathbf{U}_{\mathbf{R}}^\top [\mathbf{X}_k]_{:,a} [\mathbf{A}_k]_{:,b}^\top \mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \right) \\
&\stackrel{\text{tr}(ABC) = \text{tr}(CAB)}{=} -\text{tr} \left([\mathbf{A}_k]_{:,b}^\top \mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\tilde{\mathbf{C}}} \tilde{\mathbf{Y}}_r^\top \mathbf{U}_{\mathbf{R}}^\top [\mathbf{X}_k]_{:,a} \right) \\
&\stackrel{\text{tr}(A) = \text{tr}(A^\top)}{=} -[\mathbf{X}_k]_{:,a}^\top \mathbf{U}_{\mathbf{R}} \tilde{\mathbf{Y}}_r \mathbf{U}_{\tilde{\mathbf{C}}}^\top \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\Sigma}^\top [\mathbf{A}_k]_{:,b} \tag{1.30}
\end{aligned}$$

When stacking together the derivatives for all a and b , the gradient with respect to all entries of \mathbf{B}_k follows as

$$\begin{aligned}
& \frac{\partial}{\partial \mathbf{B}_k} \left(-\frac{1}{2} \text{vec}(\mathbf{Y}_r)^\top (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}_r) \right) \\
&= -\mathbf{X}_k^\top \mathbf{U}_{\mathbf{R}} \tilde{\mathbf{Y}}_r \mathbf{U}_{\tilde{\mathbf{C}}}^\top \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\Sigma}^\top \mathbf{A}_k \tag{1.31}
\end{aligned}$$

Closed form maximum likelihood estimates First, we rewrite the mean function by concatenating the Kronecker products of the design matrix

$$\Phi = [\mathbf{A}_1 \otimes \mathbf{X}_1, \dots, \mathbf{A}_J \otimes \mathbf{X}_J] \tag{1.32}$$

and concatenating the fixed effects

$$\boldsymbol{\beta} = \begin{pmatrix} \text{vec}(\mathbf{B}_1) \\ \dots \\ \text{vec}(\mathbf{B}_J) \end{pmatrix}, \tag{1.33}$$

where Φ is a $NP \times DM$ matrix and $\boldsymbol{\beta}$ is a $DM \times 1$ vector, with $D = \sum_j D_j$ and $M = \sum_j M_j$. It is easy to show that $\sum_{j=1}^J (\mathbf{A}_j \otimes \mathbf{X}_j) \text{vec}(\mathbf{B}_j) = \Phi \boldsymbol{\beta}$.

By setting the gradient of the log-likelihood with respect to the weight vector to zero, we can then obtain a closed form solution of the maximum-likelihood estimate $\boldsymbol{\beta}_M$:

$$\frac{\partial}{\partial \boldsymbol{\beta}_M} \left(-\frac{1}{2} (\text{vec}(\mathbf{Y}) - \Phi \boldsymbol{\beta})^\top (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} (\text{vec}(\mathbf{Y}) - \Phi \boldsymbol{\beta}) \right) = \mathbf{0} \tag{1.34}$$

$$\Phi^\top (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \Phi \boldsymbol{\beta}_M - \Phi^\top (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) = \mathbf{0} \tag{1.35}$$

$$\Phi^\top (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \Phi \boldsymbol{\beta} = \Phi^\top (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) \tag{1.36}$$

$$\boldsymbol{\beta}_M = \underbrace{\left(\Phi^\top (\mathbf{C} \otimes \mathbf{R} + \boldsymbol{\Sigma} \otimes \mathbf{I})^{-1} \Phi \right)^{-1}}_{\text{left term}} \underbrace{\Phi^\top (\mathbf{C} \otimes \mathbf{R} + \boldsymbol{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y})}_{\text{right term}} \quad (1.37)$$

We first compute the right term with respect to the j^{th} block.

$$\begin{aligned} & (\mathbf{A}_j \otimes \mathbf{X}_j)^\top (\mathbf{C} \otimes \mathbf{R} + \boldsymbol{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) \\ & \stackrel{(1.27)}{=} (\mathbf{A}_j \otimes \mathbf{X}_j)^\top (\mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R) \text{vec}(\mathbf{D} \odot \mathbf{U}_R^\top \mathbf{Y}_r \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}}) \\ & \stackrel{(1.16)}{=} (\mathbf{A}_j^\top \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{X}_j^\top \mathbf{U}_R) \text{vec}(\mathbf{D} \odot \mathbf{U}_R^\top \mathbf{Y}_r \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}}) \\ & \stackrel{(1.17)}{=} \text{vec}[\mathbf{X}_j^\top \mathbf{U}_R (\mathbf{D} \odot \mathbf{U}_R^\top \mathbf{Y}_r \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}}) \mathbf{U}_{\tilde{C}}^\top \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_\Sigma^\top \mathbf{A}_j] \end{aligned} \quad (1.38)$$

By concatenating the blocks, we obtain

$$\Phi^\top (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) = \begin{bmatrix} \text{vec}(\mathbf{X}_1^\top \mathbf{U}_R (\mathbf{D} \odot \mathbf{U}_R^\top \mathbf{Y}_r \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}}) \mathbf{U}_{\tilde{C}}^\top \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_\Sigma^\top \mathbf{A}_1) \\ \vdots \\ \text{vec}(\mathbf{X}_J^\top \mathbf{U}_R (\mathbf{D} \odot \mathbf{U}_R^\top \mathbf{Y}_r \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}}) \mathbf{U}_{\tilde{C}}^\top \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_\Sigma^\top \mathbf{A}_J) \end{bmatrix} \quad (1.39)$$

For the left term, we need to invert the block matrix $\Phi^\top (\mathbf{C} \otimes \mathbf{R} + \boldsymbol{\Sigma} \otimes \mathbf{I})^{-1} \Phi \in \mathbb{R}^{DM \times DM}$. We start with computing the $(i, j)^{\text{th}}$ block involving the i^{th} and the j^{th} fixed effects:

$$\begin{aligned} & (\mathbf{A}_i \otimes \mathbf{X}_i)^\top (\mathbf{C} \otimes \mathbf{R} + \boldsymbol{\Sigma} \otimes \mathbf{I})^{-1} (\mathbf{A}_j \otimes \mathbf{X}_j) \\ & \stackrel{(1.22)}{=} (\mathbf{A}_j \otimes \mathbf{X}_j)^\top \left[(\mathbf{U}_\Sigma \mathbf{S}_\Sigma^{\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R) (\mathbf{S}_{\tilde{C}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I}) (\mathbf{U}_{\tilde{C}}^\top \mathbf{S}_\Sigma^{\frac{1}{2}} \mathbf{U}_\Sigma^\top \otimes \mathbf{U}_R^\top) \right]^{-1} (\mathbf{A}_j \otimes \mathbf{X}_j) \\ & = (\mathbf{A}_i \otimes \mathbf{X}_i)^\top (\mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R) (\mathbf{S}_{\tilde{C}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I})^{-1} (\mathbf{U}_{\tilde{C}}^\top \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_\Sigma^\top \otimes \mathbf{U}_R^\top) (\mathbf{A}_j \otimes \mathbf{X}_j) \\ & \stackrel{(1.16)}{=} (\mathbf{A}_i^\top \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{X}_i^\top \mathbf{U}_R) (\mathbf{S}_{\tilde{C}} \otimes \mathbf{S}_R + \mathbf{I} \otimes \mathbf{I})^{-1} (\mathbf{U}_{\tilde{C}}^\top \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_\Sigma^\top \mathbf{A}_j \otimes \mathbf{U}_R^\top \mathbf{X}_j) \\ & = \sum_{c=1}^P \left(\left[\mathbf{A}_i^\top \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \right]_{:,c} \otimes \mathbf{X}_i^\top \mathbf{U}_R \right) (\mathbf{S}_{\tilde{C}}[c, c] \mathbf{S}_R + \mathbf{I})^{-1} \left(\left[\mathbf{U}_{\tilde{C}}^\top \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_\Sigma^\top \mathbf{A}_j \right]_{c,:} \otimes \mathbf{U}_R^\top \mathbf{X}_j \right) \\ & \stackrel{(1.16)}{=} \sum_{c=1}^P \left(\left[\mathbf{A}_i^\top \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \right]_{:,c} \otimes \mathbf{X}_i^\top \mathbf{U}_R \right) \left(\left[\mathbf{U}_{\tilde{C}}^\top \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_\Sigma^\top \mathbf{A}_j \right]_{c,:} \otimes (\mathbf{S}_{\tilde{C}}[c, c] \mathbf{S}_R + \mathbf{I})^{-1} \mathbf{U}_R^\top \mathbf{X}_j \right) \\ & \stackrel{(1.16)}{=} \sum_{c=1}^P \left(\left[\mathbf{A}_i^\top \mathbf{U}_\Sigma \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \right]_{:,c} \left[\mathbf{U}_{\tilde{C}}^\top \mathbf{S}_\Sigma^{-\frac{1}{2}} \mathbf{U}_\Sigma^\top \mathbf{A}_j \right]_{c,:} \otimes \mathbf{X}_i^\top \mathbf{U}_R (\mathbf{S}_{\tilde{C}}[c, c] \mathbf{S}_R + \mathbf{I})^{-1} \mathbf{U}_R^\top \mathbf{X}_j \right) \end{aligned} \quad (1.40)$$

Evaluating the term takes $O(PM^2 + PD^2N)$ time. If the number of samples is smaller than the number of traits ($N < P$), we can also explicitly sum over the samples leading to a runtime of $O(ND^2 + NM^2P)$. Inverting the left term has a runtime complexity of $O(D^3M^3)$.

1.4.4 Predictions

The mean predictor can be evaluated efficiently as follows

$$\begin{aligned}
\text{vec}(\mathbf{M}^*) &= \sum_{j=1}^J (\mathbf{A}_j \otimes \mathbf{X}_j^*) \text{vec}(\mathbf{B}_j) + (\mathbf{C} \otimes \mathbf{R}^*) (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \text{vec} \left(\mathbf{Y} - \sum_{j=1}^J (\mathbf{A}_j \otimes \mathbf{X}_j) \text{vec}(\mathbf{B}_j) \right) \\
&\stackrel{\text{def}(\mathbf{Y}_r)}{=} \sum_{j=1}^J \mathbf{X}_j^* \mathbf{B}_j \mathbf{A}_j^T + (\mathbf{C} \otimes \mathbf{R}^*) (\mathbf{C} \otimes \mathbf{R} + \mathbf{\Sigma} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}_r) \\
&\stackrel{(1.27)}{=} \sum_{j=1}^J \mathbf{X}_j^* \mathbf{B}_j \mathbf{A}_j^T + (\mathbf{C} \otimes \mathbf{R}^*) \left(\mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{U}_R \right) \text{vec}(\tilde{\mathbf{Y}}_r) \\
&\stackrel{(1.16)}{=} \sum_{j=1}^J \mathbf{X}_j^* \mathbf{B}_j \mathbf{A}_j^T + \left(\mathbf{C} \mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \otimes \mathbf{R}^* \mathbf{U}_R \right) \text{vec}(\tilde{\mathbf{Y}}_r) \\
&\stackrel{(1.17)}{=} \sum_{j=1}^J \mathbf{X}_j^* \mathbf{B}_j \mathbf{A}_j^T + \text{vec} \left[\mathbf{R}^* \mathbf{U}_R \tilde{\mathbf{Y}}_r \left(\mathbf{C} \mathbf{U}_{\Sigma} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\tilde{C}} \right)^{\top} \right] \\
&= \sum_{j=1}^J \mathbf{X}_j^* \mathbf{B}_j \mathbf{A}_j^T + \text{vec} \left(\mathbf{R}^* \mathbf{U}_R \tilde{\mathbf{Y}}_r \mathbf{U}_{\tilde{C}}^{\top} \mathbf{S}_{\Sigma}^{-\frac{1}{2}} \mathbf{U}_{\Sigma}^{\top} \mathbf{C}^{\top} \right) \tag{1.41}
\end{aligned}$$

1.5 Implementation details

OOP design LIMIX employs an object-oriented design that allows for flexibly combining different covariance functions using operations. These covariance models can be used for parameter-inference, both using gradient-based optimization and using closed-form optimization of fixed effects. The core machinery is implemented in C++ and transparent interfaces to python permit to allow construction complex genetic models at ease and without having expert knowledge.

2 Multivariate linear mixed models for statistical genetics

2.1 Introduction

In the most general formulation LIMIX models an $N \times P$ matrix \mathbf{Y} of N samples for each of P phenotypes by a multivariate linear mixed model with J fixed effects $\{\mathbf{F}_j\}$ and I random effects $\{\mathbf{U}_i\}$

$$\mathbf{Y} = \sum_{j=1}^J \mathbf{F}_j + \sum_{i=1}^I \mathbf{U}_i. \quad (2.1)$$

Each of the fixed effects and random effects is Kronecker structured. In particular, any fixed effect \mathbf{F}_i can be written as $(\mathbf{A}_i \otimes \mathbf{X}_i) \text{vec}(\mathbf{B}_i)$ where $\mathbf{A}_i \in \mathbb{R}^{P \times M}$ is the design matrix of the phenotypic effects, $\mathbf{X}_i \in \mathbb{R}^{N \times D}$ is the design matrix of the sample-specific effects, and $\mathbf{B}_i \in \mathbb{R}^{M \times D}$ is the matrix of the effect sizes; while any random effect \mathbf{U}_i is matrix-variate normal distributed, $\mathbf{U}_i \sim \mathcal{N}_{NM}(\mathbf{0}; \mathbf{R}_i, \mathbf{C}_i)$, with row covariance matrix $\mathbf{R}_i \in \mathbb{R}^{N \times N}$ and column covariance matrix $\mathbf{C}_i \in \mathbb{R}^{P \times P}$. \mathbf{R}_i and \mathbf{C}_i can be interpreted as sample-to-sample and trait-to-trait relatedness matrices due to contribution i respectively.

All the models described by equation (2.1) can be dynamically built using LIMIX and subsequently fitted to the data using the optimization framework we introduced in Section 1. In the following subsection we show how some commonly used mixed models for genetic studies can be written as in (2.1). Afterwards we review common types of genetic studies that are supported by LIMIX. In Section 2.2 we present the framework made available by LIMIX to dissect the phenotypic variability across different sources and make out-of-sample phenotype predictions. As discussed in the main text, generalization performance measured by out-of-sample prediction can be used as a model selection criterion, given the wealth of models made available by LIMIX. In Section 2.3 we describe models for genome-wide association studies (GWAS) and show how they can be extended to build a multi-locus multivariate model. Finally, in Section 2.4 we discuss the PANAMA module in LIMIX, which can be used to infer non-observed (hidden) covariates so that they can be accounted for in GWAS and variance decomposition.

2.1.1 Some recently proposed MTMM

A particular, yet multivariate case of this general model is when the same trait design is used for all fixed effects and only two random effects describing the genetic and non-genetic contributions to the phenotype are considered:

$$\mathbf{Y} = (\mathbf{A} \otimes \mathbf{X}) \text{vec}(\mathbf{B}) + \mathbf{U} + \boldsymbol{\Psi}, \quad \mathbf{U} \sim \mathcal{N}_{NM}(\mathbf{0}; \mathbf{R}, \mathbf{C}), \quad \boldsymbol{\Psi} \sim \mathcal{N}_{NM}(\mathbf{0}; \mathbf{I}, \boldsymbol{\Sigma}) \quad (2.2)$$

where \mathbf{R} is the kinship matrix, which describes the genetic relatedness between the samples, \mathbf{C} describes the genetic relatedness between the phenotypes, and $\boldsymbol{\Sigma}$ the relatedness between the phenotypes due to shared non-genetic effects, i.e. noise. Similar models have been recently employed in genetic analysis of multiple phenotypes [8, 9, 3, 7].

Considering the univariate version of the model in (2.2), the trait design matrix in the fixed effect collapses to 1 while the trait-to-trait covariance matrices reduce to single variance components. Indicating with \mathbf{y} , $\boldsymbol{\beta}$, \mathbf{u} , $\boldsymbol{\psi}$, σ_g^2 and σ_e^2 the univariate versions of \mathbf{Y} , \mathbf{B} , \mathbf{U} , $\boldsymbol{\Psi}$, \mathbf{C} and $\boldsymbol{\Sigma}$ the model reduces to

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u} + \boldsymbol{\psi}, \quad \mathbf{u} \sim \mathcal{N}(\mathbf{0}, \sigma_g^2 \mathbf{R}), \quad \boldsymbol{\Psi} \sim \mathcal{N}(\mathbf{0}, \sigma_e^2 \mathbf{I}) \quad (2.3)$$

Very similar models have been widely used in genetic studies for heritability estimation and GWAS [10, 2, 11, 12].

Instead of considering the effect from single variants as in (2.3), one can consider the effects from all SNPs in a region (e.g., a gene or a chromosome). This effect, \mathbf{r} , can be modelled as random effect whose covariance matrix \mathbf{S} is a region-based genetic relatedness matrix:

$$\mathbf{y} = \mathbf{r} + \mathbf{u} + \boldsymbol{\psi}, \quad \mathbf{r} \sim \mathcal{N}(\mathbf{0}, \sigma_r^2 \mathbf{S}), \mathbf{u} \sim \mathcal{N}(\mathbf{0}, \sigma_g^2 \mathbf{R}), \boldsymbol{\Psi} \sim \mathcal{N}(\mathbf{0}, \sigma_e^2 \mathbf{I}) \quad (2.4)$$

Similar ideas and models have been employed for general set tests [11], rare variant testing [13], and heritability estimation [14].

2.2 Variance Decomposition

In this section, we describe the variance decomposition framework made available by LIMIX that allows for dissecting the phenotypic variability across different sources for univariate (subsection 2.2.1) and multivariate analysis (subsection 2.2.2). In the same framework, LIMIX makes also available a tool for out-of-sample predictions, which we discuss in subsection 2.2.3.

2.2.1 Single-trait variance decomposition

In the model for univariate variance decomposition, fixed effects reflect covariates while random effects contributions from different sets of variants (e.g., different chromosomes or local and distal genetic contributions). The covariance matrices of these contributions are empirical covariance matrices describing genetic relatedness based on the different sets. The approach can be generalized to include non-genetic random effects whose covariance matrices are known *a priori*, allowing for correction of covariates and joint estimation of genetic and non-genetic contributions to the phenotypic variability. The model can be written as:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \sum_i \mathbf{u}_i + \boldsymbol{\psi}, \quad \mathbf{u}_i \sim \mathcal{N}(\mathbf{0}, \sigma_i^2 \mathbf{R}_i), \boldsymbol{\Psi} \sim \mathcal{N}(\mathbf{0}, \sigma_e^2 \mathbf{I}) \quad (2.5)$$

where \mathbf{R}_i is the covariance matrix for contribution i . If \mathbf{R} is a genetic contribution from set of SNPs \mathcal{S} , with a bit of abuse of notation we can write

$$\mathbf{R} = \frac{1}{C} \mathbf{W}_{:, \mathcal{S}} \mathbf{W}_{:, \mathcal{S}}^T \quad (2.6)$$

where $C = \frac{1}{N} \text{tr}(\mathbf{W}_{:, \mathcal{S}} \mathbf{W}_{:, \mathcal{S}}^T)$. The parameters of the model are the fixed effects $\boldsymbol{\beta}$ and the variance components σ_i^2 and σ_e^2 . Variance explained by the fixed effect d can be retrieved by considering $\text{var}(\mathbf{X}_{:,d} \boldsymbol{\beta}_d)$. Since one is usually interested in the relative fractions of a particular component to the total variance, variance components are normalized to sum up to 1.

2.2.2 Multiple-trait variance decomposition

Using the notation introduced in Section 2.1, the general model for multivariate variance decomposition can be written as

$$\mathbf{Y} = \sum_i (\mathbf{A}_i \otimes \mathbf{X}_i) \text{vec}(\mathbf{B}_i) + \sum_i \mathbf{U}_i + \boldsymbol{\Psi}, \quad \mathbf{U}_i \sim \mathcal{N}_{NM}(\mathbf{0}; \mathbf{R}_i, \mathbf{C}_i(\boldsymbol{\alpha}_i)), \boldsymbol{\Psi} \sim \mathcal{N}_{NM}(\mathbf{0}; \mathbf{I}, \boldsymbol{\Sigma}(\boldsymbol{\alpha}_{\Sigma})) \quad (2.7)$$

where \mathbf{U}_i indicates the genetic effect from random effect i with trait-to-trait covariance matrix \mathbf{C}_i and sample-to-sample covariance matrix \mathbf{R}_i introduced in (2.6), $\boldsymbol{\Psi}$ denotes the non-genetic contribution. We have indicated with $\boldsymbol{\alpha}_i$, $\boldsymbol{\alpha}_{\Sigma}$ the parameters of the trait-to-trait covariance matrices, which are estimated from the data. The parameters of the model are again the weight matrix of fixed effects \mathbf{B} and the variance components $\boldsymbol{\alpha} = \{\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_{\Sigma}\}$. Although LIMIX handles incomplete designs and cases where

more than 2 random effect terms are involved, in the case of 2 random effect terms and without missing values in the phenotype matrix \mathbf{Y} LIMIX exploits the Kronecker structure in the model and implements the mathematical tricks discussed in the previous chapter to conduct fast estimation of parameters.

In the most general model the trait-to-trait covariance matrices are general semi-definite positive matrices (*freeform* matrices) parameterized by $\frac{1}{2}P(P + 1)$ real values. However, such a general form may lead to overfitting as the number of parameters increases quadratically in the number of traits while the data only linearly. LIMIX circumvents this problem in two ways:

- it makes available to the user a set of different covariance matrices to perform optimization of the $P \times P$ covariance matrices;
- it allows the user to introduce a regularization on the covariance matrices

For a technical discussion about the different covariance matrices and possible regularization schemes made available by LIMIX we refer the user to the previous chapter, while here we discuss utility and interpretation of these. First of all, we highlight that using a diagonal matrix as trait-to-trait covariance matrix for random effect i is equivalent to assume that term i does not contribute to trait-to-trait correlations. A sum of a rank R matrix and spherical residual $\mathbf{C}(\boldsymbol{\theta}) = \mathbf{A}\mathbf{A}^T + \sigma^2\mathbf{I}$ (*lowrank-id*), where $\mathbf{A} \in \mathbb{R}^{P,R}$ and $\boldsymbol{\theta} = \{\text{vec}(\mathbf{A}), \sigma^2\}$, might help address overfitting problems while modelling trait-to-trait correlations as the number of parameters scales linearly with the number of traits. On the other hand, this parametrization offers limited flexibility for the diagonal elements especially for low R , yielding biased heritability estimates. A more flexible form is the sum of a low-rank term and a trait-specific independent contribution $\mathbf{C}(\boldsymbol{\theta}) = \mathbf{A}\mathbf{A}^T + \text{diag}(\mathbf{c})$ (*lowrank-diag*), where $\boldsymbol{\theta} = \{\text{vec}(\mathbf{A}), \mathbf{c}\}$ and $\mathbf{c} \in \mathbb{R}^P$. The penalization on the off-diagonal entries of the trait-trait covariance matrix, or its inverse, bridges a fully independent model (*diagonal* matrix) to the unpenalized model. To select the extent of the penalization, one can use cross validation exploiting the phenotype prediction tools discussed in the next subsection (2.2.3).

Finally, to dissect the contribution of a set of SNPs in shared and trait-specific effects, one can consider as trait-to-trait covariance matrix the sum of a block matrix, which describes pure shared effects, and a diagonal matrix, which describes pure trait-specific effects, $\mathbf{C} = a^2\mathbf{1}_{PP} + \text{diag}(\mathbf{c}^2)$. This covariance forms has been employed for the GxE variance decomposition in yeast.

2.2.3 Phenotype Predictions

The LIMIX variance decomposition model also implements a tool for phenotype predictions, which can be used to monitor overfitting and perform model selection. For example, this can help to select the set of fixed and random effects to consider, the best parametrisation of the trait-to-trait covariance matrices, or the best extent of the matrix penalization.

Following (1.14), the predictive posterior mean of the contribution all terms in model (2.7) is

$$\text{vec}(\mathbf{M}^*) = \sum_i (\mathbf{A}_i \otimes \mathbf{X}_i^*) \text{vec}(\mathbf{B}_i) + \sum_i (\mathbf{C}_i \otimes \mathbf{R}_i^*) \mathbf{K}^{-1} \text{vec} \left(\mathbf{Y} - \sum_i (\mathbf{A}_i \otimes \mathbf{X}_i) \text{vec}(\mathbf{B}_i) \right) \quad (2.8)$$

where $\mathbf{K} = \sum_i \mathbf{C}_i \otimes \mathbf{R}_i + \boldsymbol{\Sigma} \otimes \mathbf{I}$ and \mathbf{R}_i^* is the cross covariance for term i . For example, if \mathbf{R} has the form in (2.6) then we have $\mathbf{R}^* \propto \mathbf{W}_{:,S} \mathbf{W}_{:,S}^T \in \mathbb{R}^{NN^*}$ where N and N^* are the number of samples in the training and test sets respectively and \mathbf{W} and \mathbf{W}^* are their genotype.

2.3 Genome-wide Association Studies

In this section we briefly describe the methods implemented in LIMIX to perform single and multi-locus analysis in GWAS.

2.3.1 Univariate GWAS

The standard LMM considered in GWAS is

$$\mathbf{y} \sim \mathcal{N}(\mathbf{W}\boldsymbol{\alpha} + \mathbf{x}\boldsymbol{\beta}, \sigma_g^2(\mathbf{R} + \delta\mathbf{I})) \quad (2.9)$$

where \mathbf{x} is the genotypic profile of the SNP being tested, β its effect size, \mathbf{R} the sample-to-sample relatedness matrix and $\delta = \sigma_e^2/\sigma_g^2$ the signal-to-noise ratio. Testing for association of the SNP to the phenotype is done by testing $\beta \neq 0$. The δ representation of the marginal likelihood allows fast scalable methods and has been widely used in GWAS [2, 15].

The variance decomposition tool from LIMIX can be used to build more complex forms of the relatedness matrix \mathbf{R} . For example, considering the model (2.5) introduced in the previous section, we can estimate the variance components $\{\hat{\sigma}_i\}_i$ by maximum likelihood or penalized maximum likelihood and then define a new relatedness matrix $\mathbf{R} = \sum_i \hat{\sigma}_i^2 \mathbf{R}_i / \sum_i \hat{\sigma}_i^2$ that accounts for all contributions. This is equivalent to fix the contributions of each random effect term on the null model while the signal-to-noise ratio is still flexibly evaluated SNP-wise using the model in (2.9). However, building a too complex background model might explain away genetic signal leading to power reduction in GWAS if sufficient care is not taken [16]. In all our experiments we just considered the full kinship matrix as relatedness matrix, with the exception of the transcript-based eQTL analysis where we accounted for hidden confounders by using PANAMA (section 2.4).

2.3.2 Multivariate GWAS

The multivariate version of 2.9 is

$$\mathbf{Y} \sim \mathcal{N} \left((\mathbf{A}_{\text{cov}} \otimes \mathbf{W}) \text{vec}(\mathbf{B}_0) + (\mathbf{A}_1 \otimes \mathbf{x}) \text{vec}(\mathbf{B}), \sigma_g^2 (\mathbf{C} \otimes \mathbf{R} + \delta \mathbf{\Sigma} \otimes \mathbf{I}) \right) \quad (2.10)$$

where the trait-trait matrices \mathbf{C} and $\mathbf{\Sigma}$ can be estimated using the variance decomposition tool and then plugged in. This equals to estimate the variance parameters on the model with no association with the marker. However, LIMIX allows for SNP-specific estimates of the global variance of random effects σ_g^2 and the signal-to-noise ratio δ . The model in (2.10) is used to test for specific trait designs of the marker on the multivariate phenotype. To do so, the alternative model in (2.10) is compared to a null model which has the same form but characterized by a different trait design \mathbf{A}_0 for the marker effect. In the following we describe the choices of \mathbf{A}_0 and \mathbf{A}_1 to conduct the tests that are most commonly considered in multivariate GWAS and some of their extensions.

- *Any effect test*: this is a P degrees of freedom test where we test if the marker has an effect on at least one of the phenotypes ($\mathbf{A}_1 = \mathbf{I}_P$ and $\mathbf{A}_0 = \mathbf{0}_P$).
- *Common effect test*: this is a one degree of freedom test, where in the alternative model the marker has the same effect size and direction across all phenotypes ($\mathbf{A}_1 = \mathbf{1}_P$) while the null model does not contain the effect from the marker ($\mathbf{A}_0 = \mathbf{0}_P$).
- *Specific effect test*: this is a one degree of freedom test that we can use to test if the genetic marker acts specifically on the phenotype p . The design in the alternative model is a combination of a common effect and an independent effect for trait p $\mathbf{A}_1 = [\mathbf{1}_P, \mathbf{1}_{\text{trait}==p}]$, where $\mathbf{1}_{\text{trait}==p}$ return a vector with 1 at element p and 0 at all other elements, while in the null model we just have the common effect, $\mathbf{A}_0 = \mathbf{1}_P$.
- *Any specific effect test*: this is a $P - 1$ degrees of freedom test where we test whether the marker has a specific effect on at least one of the phenotypes. The alternative model describes an *any effect* from the marker ($\mathbf{A}_1 = \mathbf{I}_P$) while the null model describes a *common effect* ($\mathbf{A}_0 = \mathbf{1}_P$).

LIMIX flexibly allows considering more complex tests to really exploit the information in multivariate datasets. Here we describe an example of a more complex test. Suppose we want to jointly model T traits across E different environments, for a total of $P = TE$ trait-environment combinations. In this example, we might be interested in discovering loci that affect multiple phenotypes differently across different environment. This is a T degree of freedom test and can be performed in LIMIX by considering an alternative model where the marker has different effect sizes across traits and environments ($\mathbf{A}_1 = \mathbf{I}_P$) and a null model where the marker has same effect size across different environments for each trait ($\mathbf{A}_0 = \mathbf{I}_T \otimes \mathbf{1}_{1,E}$).

LIMIX supports also more general models for multivariate GWAS that handles non-Kronecker structured fixed effects and incomplete designs. However, such a model does not allow speed-ups and can be only used for the joint analysis of a few traits. The model implemented for such a task is practically the analogous of (2.9):

$$\mathbf{y} \sim \mathcal{N}(\mathbf{W}\boldsymbol{\alpha} + \mathbf{a}_0 \odot \mathbf{x}\beta_0 + \mathbf{a}_1 \odot \mathbf{x}\beta_1, \sigma_g^2(\mathbf{K} + \delta\mathbf{I})) \quad (2.11)$$

where \mathbf{y} is multi-phenotype vector, obtained by concatenating all observed entries of \mathbf{Y} , and an element-wise product \odot allows introducing non-Kronecker-structured designs. While the first fixed effect for SNP \mathbf{x} is contained both in the full and the null model, we test for design \mathbf{a}_1 by testing $\beta_1 \neq 0$. For example, if \mathbf{l} labels a categorical variable with values in $\{l_1, l_2, l_3\}$, setting $\mathbf{a}_0 = \mathbf{1}_P$ and $\mathbf{a}_1 = \mathbf{1}_{\mathbf{l}==l_1}$ we can test for $l_1 \times$ SNP interactions.

2.3.3 Multi-locus GWAS

All the models we introduced in this section can be extended to consider multiple loci by step-wise forward selection [12]. This is the first time multivariate analysis and step-wise forward selection are combined together. The general framework LIMIX employs multi-locus GWAS by performing the following two steps:

1. a genome-wide scan is performed using the model in (2.10)
2. if the most associated marker has P-value lower than a certain threshold the marker is added as covariate and the algorithm return to step 1, otherwise it stops here.

LIMIX allows the user to

- set the type of test to be performed in the genome-wide scans;
- set the design of the SNPs that are added as covariates;
- set a threshold either over the P-value or the Q-value of the most associated SNP as stopping criterion;
- add a maximum number of iterations as stopping criterion;
- update the parameters of the trait-to-trait covariance matrices after each inclusion using internally the variance decomposition module.

2.4 PANAMA

The PANAMA tool from LIMIX can be used to find a sample-by-sample covariance matrix that accounts for genetic and non-genetic confounders [17]. This can be done when a high number of traits are thought to share common confounding structure. This is the case of eQTL analysis, where thousands of genes most likely share sample-specific confounders. Reproducing the framework used in gaussian process latent variable models [18, 19], we assume independence of the traits (e.g., gene expression in eQTL analysis) conditions on the latent factors. After normalizing all traits to z-scores, we again consider a particular model from the general form (2.1):

$$\mathbf{y}_p = \mathbf{u} + \boldsymbol{\eta} + \boldsymbol{\psi}, \quad \forall p \quad (2.12)$$

$$\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \sigma_g^2 \mathbf{R}_g), \quad \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_c(\boldsymbol{\alpha})), \quad \boldsymbol{\psi} \sim \mathcal{N}(\mathbf{0}, \sigma_e^2 \mathbf{I})$$

where \mathbf{R}_g is the genetic relatedness matrix introduced before while $\mathbf{R}_c(\boldsymbol{\alpha})$ is the relatedness matrix due to unobserved covariates. The parameters of the model are $\{\boldsymbol{\alpha}, \sigma_g^2, \sigma_e^2\}$. As full rank matrices might overfit the data explaining away genetic signal, LIMIX models $\mathbf{R}_c(\boldsymbol{\alpha})$ as rank r matrix and allows the user to choose r . A sensible value for r can be chosen by looking at the number of PCAs of the sample-by-sample empirical covariance matrix explaining 70%-90% of the variance; see also discussion in [17].

Once \mathbf{R}_c is learned from the data, a new sample-to-sample relatedness matrix can be plugged in the GWAS tools both for single and multivariate analysis.

Bibliography

- [1] Rasmussen, C. E. & Williams, C. K. I. *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)* (The MIT Press, 2005).
- [2] Lippert, C. *et al.* FaST linear mixed models for genome-wide association studies. *Nature Methods* **8**, 833–835 (2011).
- [3] Rakitsch, B., Lippert, C., Borgwardt, K. & Stegle, O. It is all in the noise: Efficient multi-task gaussian process inference with structured residuals 1466–1474 (2013).
- [4] Liu, D. C. & Nocedal, J. On the limited memory bfgs method for large scale optimization. *Mathematical programming* **45**, 503–528 (1989).
- [5] Bernstein, D. S. *Matrix Mathematics: Theory, Facts, and Formulas* (Princeton University Press, 2009).
- [6] Petersen, K. B. & Pedersen, M. S. The matrix cookbook. Tech. Rep., Technical University of Denmark (2006).
- [7] Stegle, O., Lippert, C., Mooij, J. M., Lawrence, N. D. & Borgwardt, K. M. Efficient inference in matrix-variate gaussian models with\ iid observation noise. In *Advances in Neural Information Processing Systems*, 630–638 (2011).
- [8] Zhou, X. & Stephens, M. Efficient multivariate linear mixed model algorithms for genome-wide association studies. *Nature Methods*, *in press* (2014).
- [9] Korte, A. *et al.* A mixed-model approach for genome-wide association studies of correlated traits in structured populations. *Nature genetics* **44**, 1066–1071 (2012).
- [10] Kang, H. M. *et al.* Variance component model to account for sample structure in genome-wide association studies. *Nature genetics* **42**, 348–354 (2010).
- [11] Listgarten, J. *et al.* A powerful and efficient set test for genetic markers that handles confounding. *Bioinformatics* (2013).
- [12] Segura, V. *et al.* An efficient multi-locus mixed-model approach for genome-wide association studies in structured populations. *Nature genetics* (2012).
- [13] Wu, M. C. *et al.* Rare-variant association testing for sequencing data with the sequence kernel association test. *The American Journal of Human Genetics* **89**, 82–93 (2011).
- [14] Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for human height. *Nature genetics* **42**, 565–569 (2010).
- [15] Zhou, X. & Stephens, M. Genome-wide efficient mixed-model analysis for association studies. *Nature genetics* **44**, 821–824 (2012).
- [16] Listgarten, J. *et al.* Improved linear mixed models for genome-wide association studies. *Nature methods* **9**, 525–526 (2012).
- [17] Fusi, N., Stegle, O. & Lawrence, N. D. Joint modelling of confounding factors and prominent genetic regulators provides increased accuracy in genetical genomics studies. *PLoS computational biology* **8**, e1002330 (2012).

Bibliography

- [18] Lawrence, N. D. Gaussian process latent variable models for visualisation of high dimensional data. *Advances in neural information processing systems* **16**, 3 (2004).
- [19] Lawrence, N. D. Probabilistic non-linear principal component analysis with Gaussian process latent variable models. *The Journal of Machine Learning Research* **6**, 1783–1816 (2005).